

Predicting Manufacturing Failures with Decision Trees and SVMs

September 28, 2019

1 Predicting Manufacturing Failures with Decision Trees, Gradient Boosted Decision Trees, and Support Vector Machines

In the previous project, I generated a data flow diagram to detail the necessary steps to predict the manufacturing outcome and identify the most important features of a semiconductor manufacturing process (Figure 1). After doing some initial exploratory data analysis, it became clear that manufacturing failures are not as frequent as manufacturing successes. To overcome this, I used SMOTE to resample the training dataset and provide more failures to ensure a robust training model would be fit to the test dataset. The last project utilized LASSO for feature selection and employed a Logistic Regression model to predict manufacturing failures.

In this project, I build on the work I had done previously. The exploratory data analysis, data preparation, and SMOTE resampling were retained for context. From there I employed a few different modeling techniques such as Decision Tree, Gradient Boosted Decision Tree, and Support Vector Classifier models to determine the most optimal model to predict semiconductor manufacturing failures.

2 Import data

The data must be imported along with the labels of each column. For now I've just named the features: feature1, feature2, and so on. These data were merged with the classification column (that holds the manufacturing outcome) and the date. All NA values were set to NaN for future imputation.

```
In [1]: # import semiconductor manufacturing features and label data
import pandas as pd
url_vars = "https://archive.ics.uci.edu/ml/machine-learning-databases/secom/secom.data"
names = ["feature" + str(x) for x in range(1, 591)] # name each feature sequentially
semi_vars = pd.read_csv(url_vars, sep=" ", names=names, na_values = "NaN") # read in

url_labs = "https://archive.ics.uci.edu/ml/machine-learning-databases/secom/secom_labels"
semi_labs = pd.read_csv(url_labs, sep=" ", names = ["classification", "date"], parse_dates=

In [2]: # merge data and take a look at the first 5 rows
semi = pd.merge(semi_vars, semi_labs, left_index=True, right_index=True)
semi.head()
```

```

Out[2]:
  feature1  feature2  feature3  feature4  feature5  feature6  feature7  \
0  3030.93  2564.00  2187.7333  1411.1265    1.3602    100.0    97.6133
1  3095.78  2465.14  2230.4222  1463.6606    0.8294    100.0   102.3433
2  2932.61  2559.94  2186.4111  1698.0172    1.5102    100.0    95.4878
3  2988.72  2479.90  2199.0333   909.7926    1.3204    100.0   104.2367
4  3032.24  2502.87  2233.3667  1326.5200    1.5334    100.0   100.3967

  feature8  feature9  feature10  ...  feature583  feature584  feature585  \
0    0.1242    1.5005     0.0162  ...     0.5005     0.0118     0.0035
1    0.1247    1.4966    -0.0005  ...     0.5019     0.0223     0.0055
2    0.1241    1.4436     0.0041  ...     0.4958     0.0157     0.0039
3    0.1217    1.4882    -0.0124  ...     0.4990     0.0103     0.0025
4    0.1235    1.5031    -0.0031  ...     0.4800     0.4766     0.1045

  feature586  feature587  feature588  feature589  feature590  classification  \
0      2.3630         NaN         NaN         NaN         NaN             -1
1      4.4447      0.0096      0.0201      0.0060    208.2045             -1
2      3.1745      0.0584      0.0484      0.0148     82.8602              1
3      2.0544      0.0202      0.0149      0.0044     73.8432             -1
4     99.3032      0.0202      0.0149      0.0044     73.8432             -1

      date
0 2008-07-19 11:55:00
1 2008-07-19 12:32:00
2 2008-07-19 13:17:00
3 2008-07-19 14:43:00
4 2008-07-19 15:22:00

[5 rows x 592 columns]

```

3 Clean the dataset

After taking a look at the first 5 rows, we can see that the data have been merged properly and get an idea of how the dataset is structured. There are a steps we can take right away to clean up the dataset a bit. I started by imputing the NaNs within each feature with the median of the respective feature. Furthermore, I edited the classification column to be more straightforward by changing the name of the column to outcome and setting the successful manufacturing outcomes as 0 and the failures as a 1.

```

In [3]: # replace missing values with the median of the column
semi.fillna(semi.median(), inplace=True)
# rename the classification column to outcome in order to be more representative of manufacturing
semi.rename(columns={'classification': 'outcome'}, inplace = True)
# map the values of the outcome column to more interpretable 0 for success and 1 for failure
semi['outcome'] = semi['outcome'].map({-1: 0, 1: 1})

```

4 Exploratory data analysis

Since the dataset describes the features that contribute to whether a semiconductor is successfully manufactured or not, I first wanted to see how the outcome was distributed. There is certainly a class imbalance within the manufacturing outcomes of this dataset with only about 100 failures out of over 1500 outcomes. This will be taken care of later on with SMOTE.

To get a general sense of the features in the dataset, I obtained the summary statistics and started to dig into a few features that seemed to have an interesting summary statistics. I plotted this histograms of a few features, feature 4 and 590 appear to be have a distribution that is skewed to the right and feature 6 appears to be made up of only one value for all runs. This is interesting to note because it will most likely be removed later since it is a constant that probably doesn't affect the manufacturing outcome.

I also plotted the manufacturing outcome as a time series to determine if there was a specific time where there was any clustering of the failures around specific time frames. There are so many rows of data it's a bit hard to parse through but there does seem to be quite a few failures in August and September of 2008.

```
In [4]: from matplotlib import pyplot as plt
```

```
# plot a histogram of the manufacturing success
plt.hist(semi['outcome'])
plt.title('Distribution of Manufacturing Success and Failures')
plt.ylabel('Count')
plt.xlabel('Manufacturing Outcome')
plt.show()
```

<Figure size 640x480 with 1 Axes>

```
In [5]: # get the exact number of manufacturing success and failures
semi['outcome'].value_counts()
```

```
Out[5]: 0    1463
        1     104
        Name: outcome, dtype: int64
```

```
In [6]: # get the summary statistics for the entire dataset
semi.describe()
```

```
Out[6]:
```

	feature1	feature2	feature3	feature4	feature5	\
count	1567.000000	1567.000000	1567.000000	1567.000000	1567.000000	
mean	3014.441551	2495.866110	2200.551958	1395.383474	4.171281	
std	73.480841	80.228143	29.380973	439.837330	56.103721	
min	2743.240000	2158.750000	2060.660000	0.000000	0.681500	
25%	2966.665000	2452.885000	2181.099950	1083.885800	1.017700	
50%	3011.490000	2499.405000	2201.066700	1285.214400	1.316800	
75%	3056.540000	2538.745000	2218.055500	1590.169900	1.518800	
max	3356.350000	2846.440000	2315.266700	3715.041700	1114.536600	

	feature6	feature7	feature8	feature9	feature10	...	\
count	1567.0	1567.000000	1567.000000	1567.000000	1567.000000	...	
mean	100.0	101.116476	0.121825	1.462860	-0.000842	...	
std	0.0	6.209385	0.008936	0.073849	0.015107	...	
min	100.0	82.131100	0.000000	1.191000	-0.053400	...	
25%	100.0	97.937800	0.121100	1.411250	-0.010800	...	
50%	100.0	101.512200	0.122400	1.461600	-0.001300	...	
75%	100.0	104.530000	0.123800	1.516850	0.008400	...	
max	100.0	129.252200	0.128600	1.656400	0.074900	...	

	feature582	feature583	feature584	feature585	feature586	\
count	1567.000000	1567.000000	1567.000000	1567.000000	1567.000000	
mean	82.403069	0.500096	0.015317	0.003846	3.067628	
std	56.348694	0.003403	0.017174	0.003719	3.576899	
min	0.000000	0.477800	0.006000	0.001700	1.197500	
25%	72.288900	0.497900	0.011600	0.003100	2.306500	
50%	72.288900	0.500200	0.013800	0.003600	2.757650	
75%	72.288900	0.502350	0.016500	0.004100	3.294950	
max	737.304800	0.509800	0.476600	0.104500	99.303200	

	feature587	feature588	feature589	feature590	outcome
count	1567.000000	1567.000000	1567.000000	1567.000000	1567.000000
mean	0.021458	0.016474	0.005283	99.652345	0.066369
std	0.012354	0.008805	0.002866	93.864558	0.249005
min	-0.016900	0.003200	0.001000	0.000000	0.000000
25%	0.013450	0.010600	0.003300	44.368600	0.000000
50%	0.020500	0.014800	0.004600	71.900500	0.000000
75%	0.027600	0.020300	0.006400	114.749700	0.000000
max	0.102800	0.079900	0.028600	737.304800	1.000000

[8 rows x 591 columns]

In [7]: *# plot a histogram of the distributions of a few features*

distribution of feature4

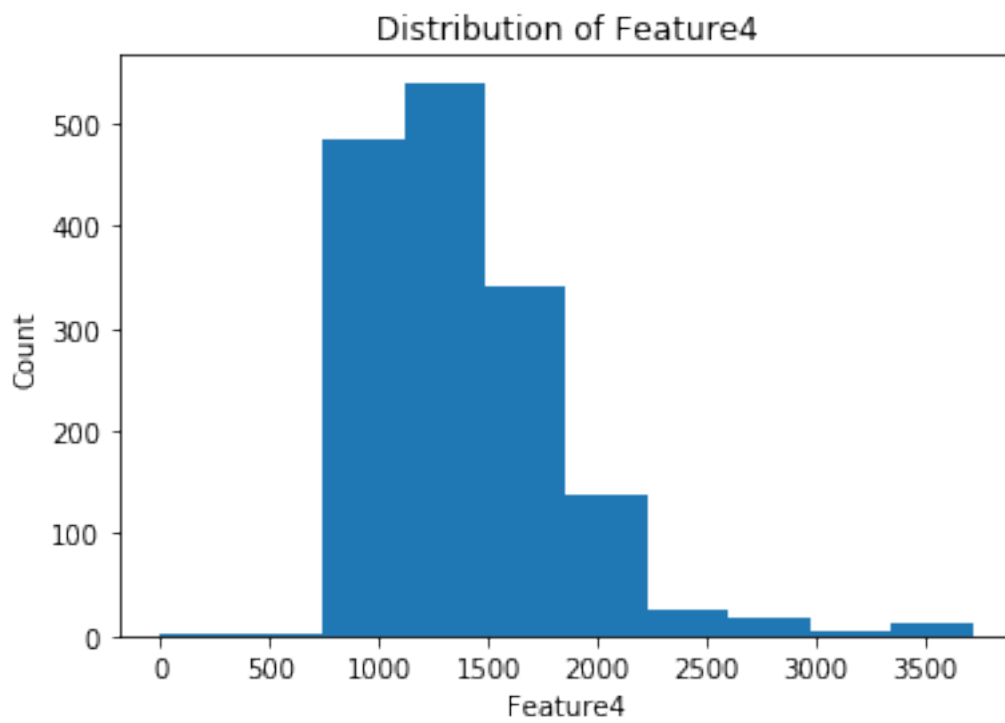
```
plt.hist(semi.feature4)
plt.title('Distribution of Feature4')
plt.ylabel('Count')
plt.xlabel('Feature4')
plt.show()
```

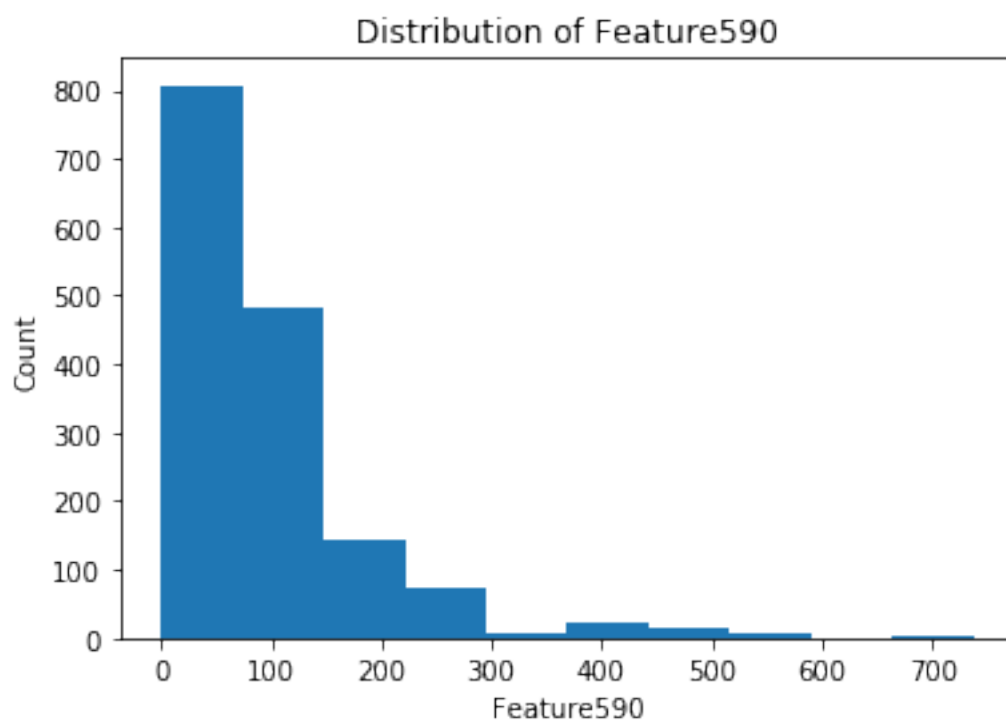
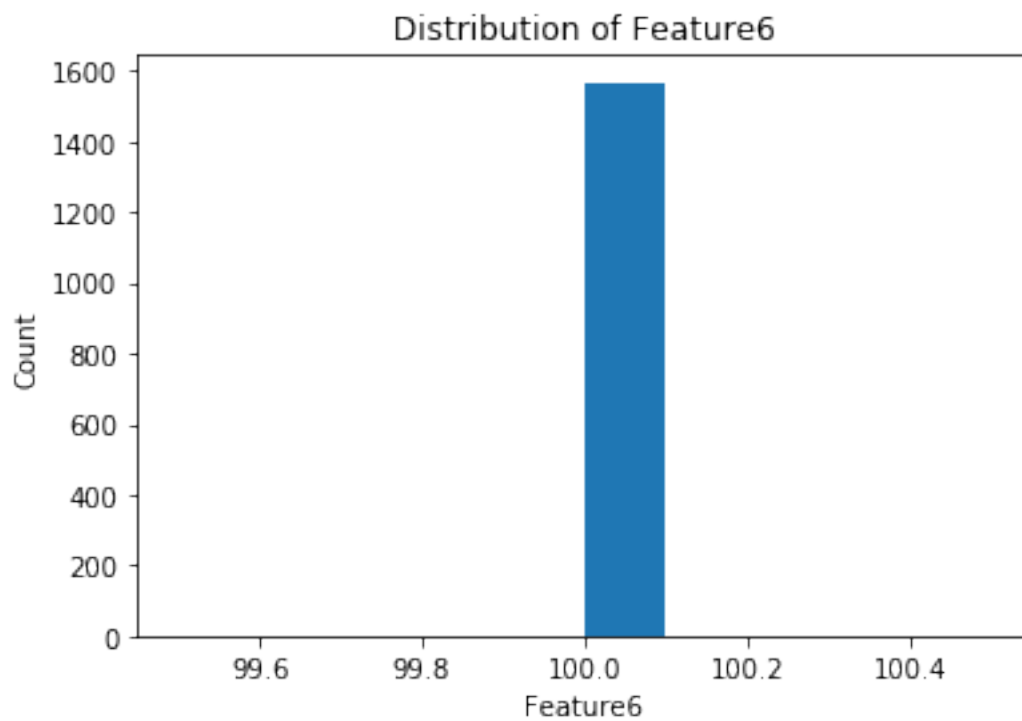
distribution of feature6

```
plt.hist(semi.feature6)
plt.title('Distribution of Feature6')
plt.ylabel('Count')
plt.xlabel('Feature6')
plt.show()
```

distribution of feature590

```
plt.hist(semi.feature590)
plt.title('Distribution of Feature590')
plt.ylabel('Count')
plt.xlabel('Feature590')
plt.show()
```





```
In [8]: # set the date column to be in the datetime format for time series analysis
        semi.loc[:, 'date'] = pd.to_datetime(semi.loc[:, 'date'])
        semi.set_index('date', inplace = True) # set index on dataset
        print(semi.head()) # print the beginning of the manufacturing dates
        print(semi.tail()) # print the end of the manufacturing dates
```

		feature1	feature2	feature3	feature4	feature5	\
date							
2008-07-19 11:55:00		3030.93	2564.00	2187.7333	1411.1265	1.3602	
2008-07-19 12:32:00		3095.78	2465.14	2230.4222	1463.6606	0.8294	
2008-07-19 13:17:00		2932.61	2559.94	2186.4111	1698.0172	1.5102	
2008-07-19 14:43:00		2988.72	2479.90	2199.0333	909.7926	1.3204	
2008-07-19 15:22:00		3032.24	2502.87	2233.3667	1326.5200	1.5334	

		feature6	feature7	feature8	feature9	feature10	...	\
date							...	
2008-07-19 11:55:00		100.0	97.6133	0.1242	1.5005	0.0162	...	
2008-07-19 12:32:00		100.0	102.3433	0.1247	1.4966	-0.0005	...	
2008-07-19 13:17:00		100.0	95.4878	0.1241	1.4436	0.0041	...	
2008-07-19 14:43:00		100.0	104.2367	0.1217	1.4882	-0.0124	...	
2008-07-19 15:22:00		100.0	100.3967	0.1235	1.5031	-0.0031	...	

		feature582	feature583	feature584	feature585	\
date						
2008-07-19 11:55:00		72.2889	0.5005	0.0118	0.0035	
2008-07-19 12:32:00		208.2045	0.5019	0.0223	0.0055	
2008-07-19 13:17:00		82.8602	0.4958	0.0157	0.0039	
2008-07-19 14:43:00		73.8432	0.4990	0.0103	0.0025	
2008-07-19 15:22:00		72.2889	0.4800	0.4766	0.1045	

		feature586	feature587	feature588	feature589	\
date						
2008-07-19 11:55:00		2.3630	0.0205	0.0148	0.0046	
2008-07-19 12:32:00		4.4447	0.0096	0.0201	0.0060	
2008-07-19 13:17:00		3.1745	0.0584	0.0484	0.0148	
2008-07-19 14:43:00		2.0544	0.0202	0.0149	0.0044	
2008-07-19 15:22:00		99.3032	0.0202	0.0149	0.0044	

		feature590	outcome
date			
2008-07-19 11:55:00		71.9005	0
2008-07-19 12:32:00		208.2045	0
2008-07-19 13:17:00		82.8602	1
2008-07-19 14:43:00		73.8432	0
2008-07-19 15:22:00		73.8432	0

[5 rows x 591 columns]

		feature1	feature2	feature3	feature4	feature5	\
--	--	----------	----------	----------	----------	----------	---

date	feature6	feature7	feature8	feature9	feature10	...	\
2008-10-16 15:13:00	2899.41	2464.36	2179.7333	3085.3781	1.4843		
2008-10-16 20:49:00	3052.31	2522.55	2198.5667	1124.6595	0.8763		
2008-10-17 05:26:00	2978.81	2379.78	2206.3000	1110.4967	0.8236		
2008-10-17 06:01:00	2894.92	2532.01	2177.0333	1183.7287	1.5726		
2008-10-17 06:07:00	2944.92	2450.76	2195.4444	2914.1792	1.5978		

date	feature582	feature583	feature584	feature585	...	\
2008-10-16 15:13:00	100.0	82.2467	0.1248	1.3424	-0.0045	...
2008-10-16 20:49:00	100.0	98.4689	0.1205	1.4333	-0.0061	...
2008-10-17 05:26:00	100.0	99.4122	0.1208	1.4616	-0.0013	...
2008-10-17 06:01:00	100.0	98.7978	0.1213	1.4622	-0.0072	...
2008-10-17 06:07:00	100.0	85.1011	0.1235	1.4616	-0.0013	...

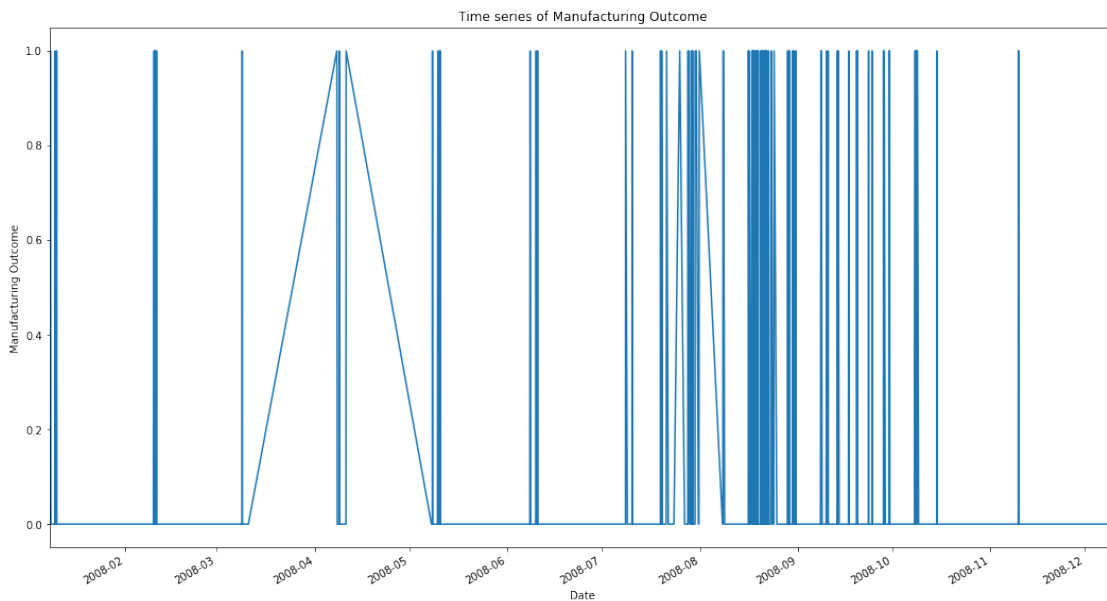
date	feature586	feature587	feature588	feature589	...	\
2008-10-16 15:13:00	203.1720	0.4988	0.0143	0.0039		
2008-10-16 20:49:00	72.2889	0.4975	0.0131	0.0036		
2008-10-17 05:26:00	43.5231	0.4987	0.0153	0.0041		
2008-10-17 06:01:00	93.4941	0.5004	0.0178	0.0038		
2008-10-17 06:07:00	137.7844	0.4987	0.0181	0.0040		

date	feature590	outcome
2008-10-16 15:13:00	203.1720	0
2008-10-16 20:49:00	203.1720	0
2008-10-17 05:26:00	43.5231	0
2008-10-17 06:01:00	93.4941	0
2008-10-17 06:07:00	137.7844	0

[5 rows x 591 columns]

```
In [9]: # plot the manufacturing outcome as a time series
ax = plt.figure(figsize=(18, 10)).gca() # define plot
semi.outcome.plot(ax = ax) # plot manufacturing outcome
ax.set_xlabel('Date')
ax.set_ylabel('Manufacturing Outcome')
ax.set_title('Time series of Manufacturing Outcome')
```


Out[9]: Text(0.5, 1.0, 'Time series of Manufacturing Outcome')



5 Prepare the data for modeling

After cleaning and getting to know the dataset a bit more, I now need to start preparing the dataset for modeling. First, I identify the outcome column as the target we are trying to predict and determine that the features will be made up of the rest of the columns in the dataset.

Next, the dataset is split up into training and test datasets further subsetting by their features and targets. The test dataset is made up of 20% of the original dataset.

```
In [10]: # establish target and features of the manufacturing data
         # set the target to the encoded manufacturing outcome column
         target = semi[['outcome']]
         # set the features as the rest of the dataset after dropping the features that are no
         feats = semi.drop(['outcome'], axis=1)
```

```
In [11]: from sklearn import model_selection

         # split original data into training and test sets
         feat_train, feat_test, target_train, target_test = model_selection.train_test_split(
                                                         test_size=0.2, random_state=6)
```

6 SMOTE

I use the SMOTE method to resample the training dataset and increase the number of failures within the dataset. This allows the model to train on more data that can help predict failures better. Because we used SMOTE, the modeling we perform will be a bit more representative of predicting manufacturing failures.

```
In [12]: import imblearn
         from imblearn.over_sampling import SMOTE
         # handle class imbalance using SMOTE
         sm = SMOTE(random_state=42)
         X_res, y_res = sm.fit_sample(feet_train, target_train)
```

```
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/utils/validation.py:761: Data
y = column_or_1d(y, warn=True)
```

7 Perform initial Logistic Regression model

I wanted to see how the model performed without handling the class imbalance or selecting specific features. I employed a Logistic Regression model on the training data and was about 95% accurate. This is a really good score but is most likely not very representative of the actual ability to detect failures due to the class imbalance. I will use the SMOTE method to resample the dataset and test the model performance again.

```
In [13]: from sklearn.linear_model import LogisticRegression

         # instantiate logistic regression model and fit to train dataset
         clf = LogisticRegression()
         logR = clf.fit(feet_train, target_train)
         # get the accuracy of the dataset
         logR.score(feet_train, target_train)
```

```
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433:
FutureWarning)
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/utils/validation.py:761: Data
y = column_or_1d(y, warn=True)
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/svm/base.py:931: ConvergenceW
"the number of iterations.", ConvergenceWarning)
```

```
Out[13]: 0.9497206703910615
```

8 Feature selection with LASSO

I used LASSO to select the most important features of the manufacturing dataset and used them to run another Logistic Regression model. LASSO selects the most important features by shrinking the coefficients of irrelevant features to 0 so they have little impact on the model. I still used the resampled training data obtained using SMOTE because it is more representative of data containing more failures. The model improved when using the selected features, as expected, to provide an accuracy of about 99%.

```
In [14]: # instantiate Logistic Regression model employing LASSO for feature selection
         clf = LogisticRegression(penalty = 'l1')
```

```

# fit model to resampled training dataset
logR = clf.fit(X_res, y_res)
# get the accuracy of the model
logR.score(X_res, y_res)

/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433:
FutureWarning)
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/svm/base.py:931: ConvergenceW
"the number of iterations.", ConvergenceWarning)

```

Out[14]: 0.9854576561163387

9 Decision Tree Model

Decision trees sort data from the "root" to the "leaves" based on an attribute splitting the data into a classification. The split attribute is systematically chosen by which ever feature provides the best decision tree model based on the coefficient selected for the algorithm.

I built a Decision Tree model with the entropy coefficient. The accuracy was in predicting manufacturing failures from this model (trained with resampled data) is about 78%. I generated a visualization of the decision tree to depict the structure of the model. Based on the visualization, it looks like the first split attribute is "feature248". The feature the model chooses to split on would be a valuable process lever to invest resources in as it is the first sign of whether a manufacturing run will fail.

While 78% is pretty good accuracy, I will try to improve the results of this model by employing an ensemble model (Gradient Boosted Decision Tree).

```

In [15]: # define decision tree parameters
nTrees = 100
max_depth = 5
min_node_size = 5
verbose = 0
learning_rate = 0.05

from sklearn.tree import DecisionTreeClassifier
# instantiate and fit decision tree classifier with the entropy coefficient
dec_tree_ent = DecisionTreeClassifier(criterion='entropy', max_depth=3)
model = dec_tree_ent.fit(X_res, y_res)

# obtain target prediction based on test features
y_predict_ent = model.predict(feat_test)

# print the accuracy score comparing predicted targets and actual test targets
from sklearn.metrics import accuracy_score
acc_ent = accuracy_score(target_test, y_predict_ent) * 100
print("Decision Tree accuracy with entropy coefficient: {}".format(acc_ent))

```

Decision Tree accuracy with entropy coefficient: 78.02547770700637%

```
In [16]: # generate visualization of decision tree
from sklearn import tree

dotfile1 = open("decisiontree_ent.dot", 'w')
tree.export_graphviz(dec_tree_ent, out_file = dotfile1, filled=True, rounded=True, fe

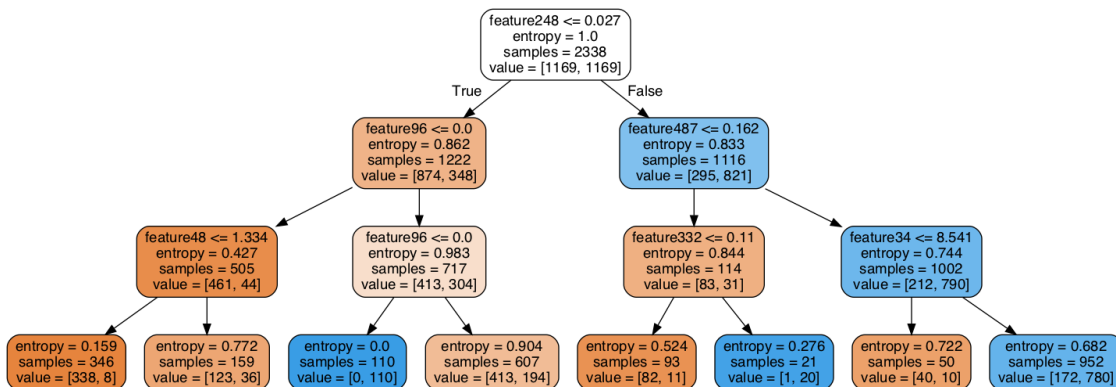
dotfile1.close()

# convert the dot file to a png
from subprocess import check_call
check_call(['dot', '-Tpng', 'decisiontree_ent.dot', '-o', 'decisiontree_ent.png'])

Out[16]: 0

In [17]: from IPython.display import Image
# print the decision tree visual
Image(filename='decisiontree_ent.png', width=500, height=500)

Out[17]:
```



10 Ensemble Model: Gradient Boosted Decision Tree

Gradient boosted decision trees combine multiple decision trees, which on their own are considered "weak learners" to generate a strong learner through multiple iterations of building decision tree models. They operate by using the entire training data to fit the residuals of the target (in this case the high risk customer) and attempt to correct the errors of the each decision tree that preceeds the current one.

I instantiated the Gradient Boosted Decision Tree model with the same parameters that I used in the original Decision Tree model above. The calculated accuracy of predicting a semiconductor manufacturing failure with the Gradient Boosted Decision Tree model is about 92%, which is quite a bit better than the Decision Tree model evaluated above (accuracy of 78%). Furthermore, I printed the importance of each feature as determined by the model. These features can be ranked to determine the process steps that require optimization or more resources.

It's encouraging to see that by adding the gradient boosting technique to leverage multiple decision trees resulted in better accuracy than just a single decision tree. The next model we will compare these to is the Support Vector Machine Classifier (SVC).

```

In [18]: # define the gradient boosted decision tree parameters
nTrees = 100
max_depth = 5
min_node_size = 5
verbose = 0
learning_rate = 0.05

from sklearn.ensemble import GradientBoostingClassifier
# instantiate and fit gradient boosted decision tree
gbm_clf = GradientBoostingClassifier(n_estimators=nTrees, loss='deviance', learning_rate=learning_rate,
                                     min_samples_leaf=min_node_size)

gbm_clf.fit(X_res, y_res)

# print feature importance
print(gbm_clf.feature_importances_)

```

[5.29174840e-03 1.74793162e-03 3.39810997e-04 4.18919388e-04
 1.35779441e-04 0.00000000e+00 4.69888394e-05 5.78149928e-06
 2.88753338e-03 3.12339221e-05 2.55613666e-04 1.52673637e-03
 2.41912217e-05 0.00000000e+00 5.72687031e-04 9.86238507e-04
 2.14771243e-04 2.67919988e-04 2.41603341e-04 1.70682688e-02
 2.08467884e-04 4.76430236e-03 3.72159066e-04 2.00230341e-03
 4.13431199e-03 6.15116657e-05 1.98741036e-04 7.63717609e-07
 2.46181079e-03 1.52253794e-03 8.69911623e-05 2.25080197e-03
 0.00000000e+00 2.46504777e-02 7.68325607e-06 7.29158020e-07
 7.38988471e-05 5.85816103e-04 0.00000000e+00 6.11657355e-03
 4.78748923e-03 4.39987323e-05 0.00000000e+00 3.52715638e-04
 1.69911297e-04 1.44271540e-03 8.99576198e-04 2.33600754e-03
 2.49677669e-05 0.00000000e+00 3.28734469e-04 3.90518189e-05
 0.00000000e+00 4.65145724e-04 5.76052976e-04 8.65867918e-03
 4.26950704e-03 3.66021520e-04 8.12239289e-04 6.93139172e-02
 1.20283381e-04 7.51638485e-04 3.32986525e-04 2.87881028e-05
 1.25352836e-02 8.77942852e-04 1.20134435e-03 9.14611889e-04
 1.20340838e-04 0.00000000e+00 1.29428042e-03 2.05811374e-03
 1.27132166e-04 8.25286728e-05 0.00000000e+00 5.27769948e-03
 1.32576105e-06 4.30763157e-03 1.82885546e-05 1.95712475e-03
 1.00738773e-04 1.28091998e-03 1.58938113e-04 2.59891997e-04
 1.86272253e-04 2.87744174e-04 8.54261056e-05 2.31636945e-04
 5.31723083e-05 1.13598347e-03 2.38975947e-03 1.82452107e-03
 2.26558275e-05 5.44708331e-04 3.67719619e-03 1.62865045e-01
 6.16168158e-06 0.00000000e+00 6.57418480e-04 2.62815366e-04
 1.08024242e-03 8.66675197e-03 1.84945036e-03 3.74710371e-04
 3.93546307e-03 4.19776195e-03 2.98359839e-04 3.72556599e-04
 8.13232327e-05 1.66648658e-04 1.06703063e-05 6.78968579e-03
 1.91464853e-02 1.93319321e-04 1.48496580e-04 9.13628868e-04
 2.40571771e-03 1.19942166e-05 9.02952892e-04 6.48192911e-04
 2.73055235e-04 6.28098398e-03 8.87622123e-04 1.44895562e-04
 5.51550131e-03 4.04388667e-03 4.79748590e-05 1.26264945e-04

0.00000000e+00 5.95264678e-03 3.51172368e-03 7.77798443e-03
1.03265001e-04 1.17139966e-03 2.06044315e-04 5.71321511e-05
7.98588686e-04 2.16610599e-04 3.81710952e-04 4.95456176e-05
6.01088598e-05 0.00000000e+00 4.21937783e-04 2.12376345e-05
2.90884790e-04 1.28261106e-03 4.46338841e-03 1.39550465e-04
1.99521568e-05 0.00000000e+00 2.31542831e-03 2.61770792e-05
2.15352635e-03 1.06162911e-03 6.85837895e-05 2.81333154e-05
7.66869667e-04 0.00000000e+00 5.70016480e-05 9.29070524e-04
1.11687108e-03 8.65437658e-06 3.68293306e-04 8.85809448e-04
0.00000000e+00 7.51845717e-05 1.47011692e-05 2.49446329e-04
2.58429459e-04 5.75043170e-05 4.19826647e-04 8.93755455e-05
0.00000000e+00 1.64194111e-04 0.00000000e+00 7.69927712e-04
2.34395174e-04 7.15371822e-04 0.00000000e+00 0.00000000e+00
1.08703794e-04 1.10575302e-05 3.40421884e-05 1.31206987e-03
3.59415352e-05 4.27861170e-04 0.00000000e+00 1.19736740e-05
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 7.91328088e-04
3.96458305e-04 2.14051677e-04 2.22276237e-05 3.30243014e-03
2.60353864e-04 6.52123284e-04 5.11890990e-04 4.33944119e-05
7.76399506e-06 7.36267552e-04 0.00000000e+00 6.05193989e-05
0.00000000e+00 0.00000000e+00 5.63796519e-04 0.00000000e+00
7.63638626e-05 1.31062568e-03 3.32732996e-07 8.18554983e-05
8.95962627e-04 6.29122732e-04 1.53388684e-03 6.44625343e-04
0.00000000e+00 2.73752124e-05 1.11014036e-04 2.10353893e-03
7.12983963e-05 1.69860007e-04 0.00000000e+00 4.53673751e-04
1.75211385e-05 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 3.35646927e-04 2.53266892e-04
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
7.46535481e-05 1.04316073e-04 0.00000000e+00 2.34667963e-02
1.12387430e-06 0.00000000e+00 1.16319769e-04 1.34122737e-04
0.00000000e+00 6.05560457e-03 4.04452156e-05 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
1.25516875e-04 2.64047254e-04 2.61878316e-05 1.29449780e-04
8.34130475e-04 2.46768479e-05 7.93530707e-04 7.37360601e-05
0.00000000e+00 3.70189718e-05 1.71560093e-03 9.48260322e-06
1.66674493e-03 4.25036357e-04 1.15356704e-03 3.31206330e-05
0.00000000e+00 6.30880725e-05 3.82710641e-04 1.15125795e-02
0.00000000e+00 0.00000000e+00 3.09934585e-04 4.33617064e-05
9.38051417e-04 3.83061555e-04 2.91480925e-05 1.34809888e-03
1.08365614e-06 8.47799893e-07 2.71182433e-03 1.25207174e-03
9.34591571e-05 1.40602399e-03 1.62384427e-03 2.06859616e-04
6.38059234e-04 1.02279356e-04 1.52857933e-04 3.74253618e-03
1.57312063e-03 4.41124180e-03 2.26123278e-03 5.46122353e-05
1.93442835e-03 0.00000000e+00 0.00000000e+00 0.00000000e+00
2.44727770e-05 0.00000000e+00 0.00000000e+00 5.93039849e-04

2.87102870e-04 3.32645718e-06 0.00000000e+00 2.43992452e-04
 2.67514364e-04 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 8.29089933e-03
 1.76334813e-03 5.76418193e-04 6.88397420e-06 0.00000000e+00
 3.88194724e-04 7.84678652e-03 4.07116667e-05 1.14157350e-04
 5.16850758e-05 1.66581041e-03 0.00000000e+00 1.19932487e-03
 6.85497556e-04 7.24229241e-03 2.24798489e-02 0.00000000e+00
 1.69376776e-03 1.02046102e-03 2.49173290e-03 1.13061351e-04
 7.68458220e-07 2.35672146e-04 8.34311204e-08 7.29238003e-05
 5.04454853e-05 0.00000000e+00 4.42365985e-04 0.00000000e+00
 2.26945660e-04 2.27822382e-05 1.08194650e-05 4.50800607e-04
 0.00000000e+00 3.74538172e-03 2.57733624e-05 9.82034928e-04
 4.55442840e-05 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
 3.62474552e-03 3.85668823e-04 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 9.74463661e-04 2.40688607e-02
 1.43335934e-03 5.16843235e-03 1.47027352e-04 0.00000000e+00
 5.31858196e-04 3.61064719e-04 7.37639485e-03 2.19591289e-03
 4.63049505e-05 3.06871579e-04 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 2.14477633e-04 2.65289247e-03 1.11110058e-03
 2.95917340e-06 6.09154427e-05 4.68003971e-04 4.42455138e-03
 5.69452396e-04 1.40589860e-03 0.00000000e+00 1.23217952e-04
 7.38523518e-04 0.00000000e+00 3.62771301e-03 3.23639015e-04
 1.63335009e-04 0.00000000e+00 0.00000000e+00 1.89879376e-03
 5.30652880e-05 1.80614917e-04 8.88432380e-05 0.00000000e+00
 2.44716432e-04 8.07518577e-04 4.80963958e-04 9.06919644e-04
 3.81004491e-04 6.31249955e-03 1.31833518e-03 2.03968228e-04
 1.27970638e-05 1.86466953e-04 1.96086102e-05 5.29318067e-04
 7.24698929e-05 2.87393734e-08 0.00000000e+00 8.83247844e-06
 1.01150270e-05 1.57672062e-04 1.14871372e-03 1.78686092e-05
 0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
 4.55932394e-05 1.82954337e-05 5.60225849e-04 4.97799167e-03
 3.18571959e-05 2.19027686e-05 0.00000000e+00 3.25882529e-04
 3.29856977e-04 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 1.48084333e-03
 5.21722369e-03 8.37116765e-04 0.00000000e+00 1.69579515e-04
 1.36004743e-03 2.84500412e-03 1.54218034e-05 4.80679230e-04
 2.94364223e-04 3.61159664e-03 0.00000000e+00 1.27638047e-05
 7.32921962e-04 0.00000000e+00 9.47375103e-04 6.51653312e-04
 7.42925912e-04 2.08996391e-03 3.33839655e-02 8.50781268e-03
 4.64934071e-04 7.51192111e-03 3.44958388e-05 4.38171564e-04
 4.85872635e-05 1.05496022e-05 1.67562903e-04 7.48231243e-07
 8.79749001e-04 5.62825352e-05 0.00000000e+00 1.47733391e-04
 3.38400462e-03 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
 0.00000000e+00 0.00000000e+00 5.57662143e-03 2.07671238e-02

```

0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
2.65284236e-04 1.34461619e-04 1.54476969e-03 1.18728494e-01
0.00000000e+00 0.00000000e+00 5.53395873e-04 3.25601648e-04
9.83245966e-04 1.42687355e-04 0.00000000e+00 9.39467047e-05
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
6.41159263e-07 8.88739223e-03 7.73282214e-04 6.07607461e-05
0.00000000e+00 3.20579969e-04 0.00000000e+00 1.62601885e-03
4.58821700e-04 8.99681084e-04 0.00000000e+00 1.20453097e-03
9.70295991e-04 4.36990590e-05 5.31030462e-05 0.00000000e+00
1.69767735e-05 1.48362309e-05 2.50013559e-04 2.25354740e-04
5.22840538e-03 6.28919044e-04 1.04149381e-03 8.36133446e-03
1.56953829e-04 3.48764067e-04 2.16450753e-04 3.60090528e-05
1.09950434e-03 9.72274452e-05 8.60293653e-07 1.64511413e-03
4.01754113e-04 3.51697527e-04 2.86501336e-04 4.67689038e-04
1.39671447e-03 5.64599980e-05 4.29189185e-05 0.00000000e+00
1.18062824e-05 2.74924512e-05 3.49503400e-05 2.92258467e-05
4.98378035e-05 2.52754569e-04 1.07443657e-04 0.00000000e+00
2.14748507e-06 4.42973051e-04]

```

```

In [19]: # obtain gradient boosted decision tree target prediction based on test features
gbm_pred = gbm_clf.predict(feat_test)

# print accuracy by comparing the predicted targets and the actual targets from the t
gbm_score = accuracy_score(target_test, gbm_pred)
gbm_accuracy = gbm_score*100
print("Gradient Boosted Decision Tree Classifier sccuracy: %.1f%%"%gbm_accuracy)

```

Gradient Boosted Decision Tree Classifier sccuracy: 91.1%

11 Support Vector Classification (SVC) Model

After investigating the performance of a simple decision tree model and improving it's accuracy through gradient boosting ensemble modeling, I take a look at how well a Support Vector Classifier predicts semiconductor manufacturing failures.

The Support Vector Classifier model behaves a bit differently than a decision tree in that a separating hyperplane is generated to define the two classes. The labeled training data (resampled using SMOTE in this case) is fed into the model and an the algorithm outputs the optimal hyperplane categorizing the two classes.

In order to obtain the most optimal hyperplane, there are a few hyper parameters we need to tune. These include the kernel, cost, and gamma. The kernel defines how the data are transformed to be respresented as a vector to then be divided by a hyper plane. The cost represents the level of penalty given to the error term, in this case I selected 0.9 which is quite stringent. This means that I'm imposing a higher cost for misclassifications and employing what is called a "hard margin". If the cost value was lower, it would allow the model more leeway and provide a "soft margin". Soft

margins tend to be more general and have a lower sensitivity for noise. Finally, the gamma parameter defines how influential a single training example is, low gamma meaning a single training example is 'far' and high values meaning they are 'close'. For now, I set gamma to be equal to 5.

To predict semiconductor manufacturing failures, I set a hard margin (cost = .9) to determine what the accuracy of the model is when there is a large amount of penalty to the error term. I tried out a few different kernels to observe how they affect the accuracy of the model. I started with a linear kernel just to get an initial read on how the model was performing. Using the classification report, we can evaluate how the model performed by looking at the precision, recall, and f1-score. I will be primarily looking at the f1-score since it combines precision and recall. Right off the bat we can see that the model performs quite well with an f1-score of 0.9. I'll try to improve this score by tweaking the hyperparameters and trying other kernels.

```
In [20]: # take a best guess at the hyper parameters to use
cost = .9 # penalty parameter of the error term
gamma = 5 # defines the influence of input vectors on the margins
```

```
In [21]: from sklearn import svm, metrics
from sklearn.metrics import classification_report

# test a LinearSVC with the initial hyper parameters
clf1 = svm.LinearSVC(C=cost).fit(X_res, y_res)
clf1.predict(feet_test)
print("LinearSVC")
print(classification_report(clf1.predict(feet_test), target_test))
```

```
LinearSVC
```

	precision	recall	f1-score	support
0	0.99	0.94	0.96	312
1	0.00	0.00	0.00	2
micro avg	0.93	0.93	0.93	314
macro avg	0.50	0.47	0.48	314
weighted avg	0.99	0.93	0.96	314

```
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/svm/base.py:931: ConvergenceWarning:
  "the number of iterations.", ConvergenceWarning)
```

12 SVC: rbf and poly Kernels

There is a "kernel trick" you can employ when fine tuning SVCs. Ideally you want to use a linear kernel when you data is linear. However, if you have data that is too non-linear, you can transform your input variables so that the shape of your dataset becomes more linear. To do this you can select different kernels to transform your dataset.

In this case, I tried the rbf and poly kernels to see if they perform better than the linear kernel above. It looks like the rbf kernel (f1-score of 0.97) performs better than the linear kernel (0.90) but the poly kernel actually performs worse (f1-score of 0.86). I'll move forward with the rbf kernel and see if adjust the cost parameter affects the model performance.

```
In [22]: # test linear, rbf and poly kernels
        for k in ('rbf', 'poly'):
            clf = svm.SVC(gamma=gamma, kernel=k, C=cost).fit(X_res, y_res)
            clf.predict(feet_test)
            print(k)
            print(classification_report(clf.predict(feet_test), target_test))
```

rbf

```
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/metrics/classification.py:114:
'recall', 'true', average, warn_for)
```

	precision	recall	f1-score	support
0	1.00	0.94	0.97	314
1	0.00	0.00	0.00	0
micro avg	0.94	0.94	0.94	314
macro avg	0.50	0.47	0.48	314
weighted avg	1.00	0.94	0.97	314

poly

	precision	recall	f1-score	support
0	0.90	0.94	0.92	280
1	0.20	0.12	0.15	34
micro avg	0.85	0.85	0.85	314
macro avg	0.55	0.53	0.53	314
weighted avg	0.82	0.85	0.84	314

13 SVC: Adjusting Cost Parameter

I decided to move forward with the rbf kernel to try and tweak some of the other hyper parameters in hopes of improving the model further. I updated the cost hyper parameter to be a bit more lenient (0.6) to get an idea of how much of a trade off we can expect if we are more lenient in our predictions. It's unlikely that the manufacturing plant would like to have a softer margin in predicting semiconductor failures if it doesn't dramatically improve the performance of the

model. We currently have a pretty solid model with an f1-score of 0.97 but we are doing our due diligence and check if adjusting the cost parameter has any impact.

After running the model again with the rbf kernel and a cost of 0.6, the f1-score does not change at all (0.97). Moving forward, it would be best to impose a hard margin and keep the cost parameter at 0.9 and use the rbf kernel in the SVC to predict semiconductor manufacturing failures.

```
In [23]: # rbf SVC with a lower cost parameter
rbf = svm.SVC(gamma=gamma, kernel='rbf', C=0.6).fit(X_res, y_res)
rbf.predict(feet_test)
print("rbf")
print(classification_report(rbf.predict(feet_test), target_test))
```

```
rbf
```

	precision	recall	f1-score	support
0	1.00	0.94	0.97	314
1	0.00	0.00	0.00	0
micro avg	0.94	0.94	0.94	314
macro avg	0.50	0.47	0.48	314
weighted avg	1.00	0.94	0.97	314

```
/Users/caseythayer/anaconda3/lib/python3.6/site-packages/sklearn/metrics/classification.py:114:
'recall', 'true', average, warn_for)
```

14 Conclusions

I built off of the previous project to employ a few additional models to predict semiconductor manufacturing failures. In the previous experiment I designed a data flow diagram to clearly lay out the steps in predicting manufacturing failures using machine learning. I read in the data and did some exploratory data analysis. At this point, it was evident that the number of failures in the dataset were minimal and to improve the robustness of the training dataset, SMOTE resampling was performed to increase the number of failures in the training dataset. The previous project ran a Logistic Regression model and achieved an accuracy of 95%. LASSO feature selection identified the most important features of the manufacturing process so that the company could invest more resources into improving specific process levers and ultimately prevent manufacturing failures.

In this project, I retained the data flow diagram, cleaning / exploration of the data, and the SMOTE resampling of the training data to be fed into Decision Tree, Gradient Boosted Decision Tree, and Support Vector Classifier models. The purpose of this project is to evaluate these three models and identify the model that is best suited to predict semiconductor manufacturing failures.

Decision Tree Conclusions * The decision tree model splits the data based on specific attributes to determine a final classification. In this case, predicting semiconductor manufacturing failures, I built a decision tree model with the entropy coefficient and trained the model on the resampled training dataset. * The decision tree model accuracy in predicting semiconductor manufacturing

failures was 78%. This is a decent model but still did not perform as well as the Logistic Regression model employed in the previous project (accuracy of 95%). * I visualized the decision tree to observe which attributes the algorithm was splitting the data from. These features could be important process levers to invest resources in tightening up since according to the decision tree their behavior directly determines whether the model predicts a manufacturing success or failure (the top three features are: feature248, feature96, and feature487). * In an effort to improve the performance of the decision tree model I decided to try an ensemble model called Gradient Boosted Decision Tree.

Gradient Boosted Decision Tree Conclusions: * Gradient Boosted Decision Trees take a simple learner (decision tree model) and combine a number of them to improve the robustness in predicting a specific outcome. * In this case, the gradient boosted decision tree utilized a number of decision trees to predict semiconductor manufacturing failures (trained on the resampled training dataset). The accuracy was quite an improvement from the Decision Tree model described above, with an accuracy of about 92%. The performance of this model is closer to the accuracy of 95% observed in the Logistic Regression model performed in the previous project.

Support Vector Classifier Conclusions: * Support Vector Classifiers take a slightly different approach than decision trees in that they identify an optimal hyper plane to classify a specific outcome. To determine the optimal hyper plane, it's important to tune the hyper parameters of the model (kernel, cost, and gamma). * In this case, I tried out a few different kernels and adjusted the cost hyper parameter to build the a SVC to predict semiconductor manufacturing failures. To evaluate the model, I used the f1-score to determine if more tweaking was required. * First, I set my cost parameter to 0.9 (a hard margin) and gamma to 5. I tried a linear kernel and obtained a pretty good f1-score of 0.9. This is on par with what we saw in the Gradient Boosted Decision Tree and Logistic Regression models. * Next, I tried to improve the f1-score by trying out the rbf and poly kernels. These kernels with transform the input data to build a more optimal hyper plane. The rbf kernel performed better than the linear kernel with an f1-score of 0.97. However, the poly kernel performed worse with an f1-score of 0.86. * Finally, I tried to loosen the cost parameter using the rbf kernel, which results in less penalty to the error term for incorrect predictions and the f1-score remained the same. * Therefore, after tuning the hyper parameters and selecting a kernel, an SVC model that utilizes an rbf kernel, with a cost of 0.9 will provide very strong results in predicting manufacturing failures with an f1-score of 0.97.

Final Conclusions: * The most optimal model throughout both projects is the SVC model with the rbf kernel and a hard margin (cost = 0.9). The f1-score provides a good representation of both precision and recall of a model and a score of 0.97 is very good. The Logistic Regression model from the previous project still performed very well with an accuracy of 95% and could be used as well since it's a bit more straightforward of a model. * The decision tree models provided additional information by visualizing which attributes the model splits with and the feature importance obtained by the Gradient Boosted Decision Tree model * Overall, through a combination of these analyses a company can derive the knowledge necessary to identify process levers that require additional optimization or resourcing and continue to track manufacturing failures and performance through the implementation of these modeling techniques.